

# Visual SLAM-based Localization and Navigation for Service Robots: The Pepper Case

Cristopher Gómez<sup>1\*\*</sup>, Matías Mattamala<sup>1\*\*</sup>, Tim Resink<sup>3\*\*</sup>, and Javier Ruiz-del-Solar<sup>1,2</sup>

<sup>1</sup> Department of Electrical Engineering, Universidad de Chile, Chile  
`{cristopher.gomez, mmattamala, jruizd}@ing.uchile.cl`

<sup>2</sup> Advanced Mining Technology Center, Universidad de Chile, Chile

<sup>3</sup> Delft University of Technology, The Netherlands  
`p.w.resink@student.tudelft.nl`

**Abstract.** We propose a Visual-SLAM based localization and navigation system for service robots. Our system is built on top of the ORB-SLAM monocular system but extended by the inclusion of wheel odometry in the estimation procedures. As a case study, the proposed system is validated using the Pepper robot, whose short-range LIDARs and RGB-D camera do not allow the robot to self-localize in large environments. The localization system is tested in navigation tasks using Pepper in two different environments: a medium-size laboratory, and a large-size hall.

## 1 Introduction

Pepper is the official robot used in the RoboCup@Home Standard Platform League. It presents several advantages for human-robot interaction such as its friendly appearance but has important limitations such as its reduced sensing and computing capabilities. In contrast to custom robots which generally rely on expensive LIDARs for metric localization and navigation, which work in both indoor and outdoor environments, Pepper has short-range LIDARs and an RGB-D camera that provide reliable localization only in small indoor rooms, being unable to provide useful information to localize the robot in large environments. This is a big deal for Pepper, which is expected to be used not only in homes, but also in public places like hospitals, shopping malls, and schools.

To address this issue, we built upon the recent advances of visual SLAM systems to develop a visual-SLAM based self-localization solution aided by wheel odometry, which allows Pepper to self-localize and navigate in large environments. The reason to include odometry in the visual estimation procedures is to recover the metric scale (unknown in typical pure-visual schemes) and to make the visual system more robust to tracking failures. This is vital for navigation tasks that require a “continuous” localization hypothesis to work. The proposed solution is based on an open-source visual SLAM system, ORB-SLAM [10], which

---

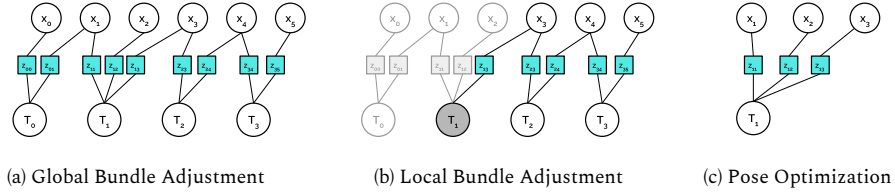
<sup>\*\*</sup> Equal contribution

is extended by the inclusion of the wheel's odometry in the estimation procedures.

In Section 2 we present a brief overview of modern SLAM systems. Then, in Section 3 we describe some basic notation as well as relevant characteristics of the Pepper robot. Afterwards, we present our localization and navigation approach in Section 4. In Section 5 we present two experiments of localization and navigation with the Pepper robot in different environments. Finally, Section 6 concludes the work with discussion and recommendations for future developments along this line.

## 2 Visual SLAM

Visual SLAM has been a hot topic during the last years since it presents a low-cost solution for applications that require localization and mapping features such as augmented reality, virtual reality, and autonomous systems (e.g. autonomous cars, inspection drones). Being originally formulated as a filtering problem, nowadays optimization-based approaches are preferable by its superior accuracy at similar computational cost [13]. Optimization-based approaches model the problem as a *factor graph* which probabilistically relates several variables -such as poses and landmarks-, by the so-called *factors*, that correspond to sensor measurements or physical constraints between the variables [1]. An example of a visual SLAM system is shown in Figure 1.



**Fig. 1.** Different factor graphs related to optimization approaches in ORB-SLAM. Circles denotes variables such as map points and keyframes within a visual SLAM scheme; white are active, gray fixed. Squares denote factors or measurements.

The factor graph can be formulated as a non-linear least squares problem [1] that aims to find the states  $\mathcal{X}^* = \mathcal{X}_1, \dots, \mathcal{X}_m$  that minimize the error between the measurements  $\mathcal{Z}_i$ , and an *observational model*  $h_i(\mathcal{X}_i)$  that "predicts" the expected measurement given the state  $\mathcal{X}_i$ <sup>4</sup>:

$$\mathcal{X}^* = \arg \min_{\mathcal{X}} \sum_{i=1}^m \|h_i(\mathcal{X}_i) \boxminus \mathcal{Z}_i\|_{\Omega_i}^2 \quad (1)$$

<sup>4</sup> The operator  $\boxminus$  generalizes the concept of subtraction for non-Euclidean spaces. Please refer to Hertzberg *et al.* [5] for a complete treatment.

The same formulation holds for the visual case, where the states correspond to selected camera poses of the trajectory -keyframes- and also the map representation -3D points, surfels, voxels, etc-, and the measurements are reprojections of the map into the image plane.

Regarding some actual systems, different solutions have been developed for monocular and stereo/depth sensors. We are concerned about monocular solutions since cameras are ubiquitous in current robots while being “cheap” sensors in comparison to the other two; they also can work in both indoor and outdoor environments. Monocular visual SLAM systems are either *feature-based* that utilize just some features in the image, such as ORB-SLAM [10], or direct methods that exploit the complete information from every image as in LSD-SLAM [3].

The main issue with monocular systems is that they require a moving camera in order to estimate the depth of the scene, as well as depending on an unknown scale factor that maps the estimated states to physically consistent dimensions. The typical approach to solve the problem relies on the usage of different sources of information that provides the scale, such as inertial measurements units (IMU); however, this increases the computational requirements of the estimation problem, since the number of states increases [12].

The utilization of visual localization systems in the RoboCup@Home has been disregarded since most of the custom robots could afford accurate but expensive LIDAR systems [7] [2], which provide a simpler solution. Nevertheless, since the range of Pepper’s LIDARs and depth camera are defined by the manufacturer, and the RoboCup@Home SSPL (Social Standard Platform League) forbids the use of additional sensors, it is unfeasible for the robot neither localize nor navigate in large environments. For this reason, we propose a visual approach for the localization problem based on an open-source visual SLAM system, ORB-SLAM [10], and we present a strategy to solve visual SLAM issues (mainly the lack of a metric scale) by aiding the system with wheel odometry measurements.

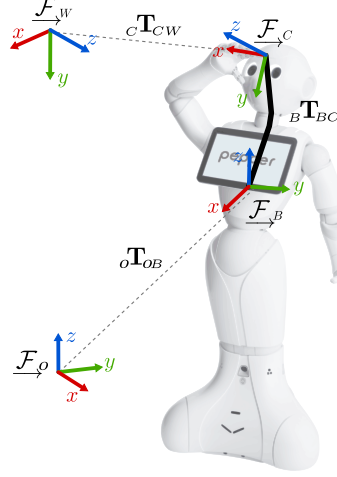
### 3 Platform, coordinate systems and notation

#### 3.1 Notation

To prevent confusion in notation, we follow the conventions of Paul Furgale<sup>5</sup>:

- Coordinate frame A is notated as  $\underline{\mathcal{F}}_A$ .
- Homogeneous transformation matrix  ${}_o\mathbf{T}_{wC} \in \text{SE}(3)$  represents the pose of the camera frame  $\underline{\mathcal{F}}_C$  with respect to the world frame,  $\underline{\mathcal{F}}_W$ , seen from frame  $\underline{\mathcal{F}}_O$ . A vector expressed in world frame W,  ${}_W\mathbf{v}$  can be hereby transformed to the camera frame C by the rotation matrix  $\mathbf{R}_{wC} \in \text{SO}(3)$ , as  ${}_C\mathbf{v} = \mathbf{R}_{wC} {}_W\mathbf{v}$
- The homogeneous transformation matrix  ${}_C\mathbf{T}_{wC}$  will be abbreviated to  $\mathbf{T}_{wC}$  for reader convenience unless otherwise indicated.

<sup>5</sup> <http://paulfurgale.info/news/2014/6/9/representing-robot-pose-the-good-the-bad-and-the-ugly>



**Fig. 2.** Coordinate frames used in this work, we follow the classic conventions with X red, Y green and Z blue.  $\mathcal{F}_c$  denotes the camera frame,  $\mathcal{F}_o$  the odometry frame and  $\mathcal{F}_B$  the body's. Pepper picture is based on Philippe Dureuil's.

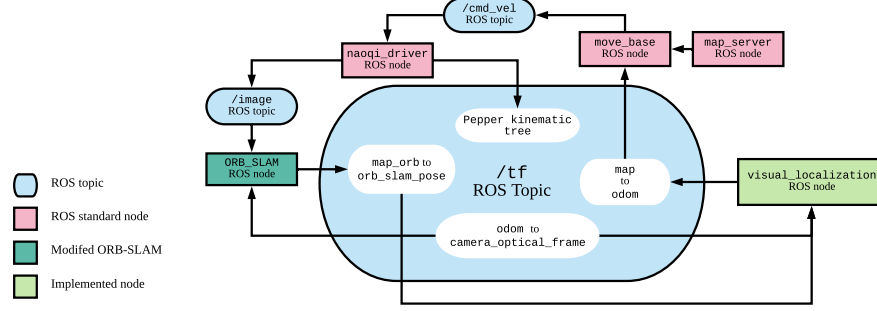
### 3.2 Pepper robot

Pepper is a wheeled humanoid platform. It has a mobile omnidirectional base and 20 degrees of freedom, including an actuated pelvis and neck. It has two Omnivision OV5640 cameras, located in the forehead and the mouth, in addition to an RGB-D sensor in the eyes. Additionally, the base has three LIDARs and an IMU. In order to access the sensors and control the robot, Softbank provides an API to its middleware, NAOQi.

Since we base our system in the ROS framework, we access sensing and perform control through the `naoqi_driver` ROS package. We principally use the images from the forehead camera at a 640x480 pixels resolution, as well as the internally computed odometry measurements; the algorithmic details about the latter are unknown to the user.

We considered two main reference frames for this work (Figure 2): On the one hand, the odometry frame denoted by  $\mathcal{F}_o$  describes the pose of the robot relative to the initial pose, as defined in ROS REP 105<sup>6</sup>. We use this frame to describe the pose of the robot's torso (body) computed by the internal wheel odometry, denoted by  ${}_o\mathbf{T}_{oB}$ . On the other, ORB-SLAM has its own reference frame (world)  $\mathcal{F}_w$  that depends on the initialization of the system, hence it may change every time the system is reset. The estimate provided by ORB-SLAM corresponds to the world position with respect to the camera pose,  ${}_c\mathbf{T}_{cw}$ .

<sup>6</sup> <http://www.ros.org/reps/rep-0105.html>



**Fig. 3.** Overview of our proposed system. The camera images are feed to the ORB-SLAM system together with the camera position with respect to the odometry frame (*odom*). An estimated camera position with respect to an arbitrary fixed frame is given as output by ORB-SLAM. The `visual_localization` node takes this information and the Pepper kinematic information to compute a transformation between the standard fixed frame *map* and the *odom* frame.

## 4 Localization and navigation system

Our visual SLAM-based localization and navigation system for Pepper consist of three main components, which are shown in Figure 3. Firstly, an ORB-SLAM-based localization and mapping system, which uses a single RGB camera located in Pepper’s forehead, and odometry measurements computed by the proprietary Pepper’s software. The second component correspond to the `visual_localization`<sup>7</sup> ROS node that transforms ORB-SLAM’s camera pose estimate to a transformation between the standard *map* frame and the *odom* frame. Finally, the node `move_base`<sup>8</sup> executes the navigation process.

### 4.1 ORB-SLAM-based localization

Our localization system maintains the same software architecture with 3 parallel threads, original from ORB-SLAM2 [11]: incoming images are processed in the *Tracking* thread, creating new map points and estimating the current camera pose  $\mathbf{T}_{cw}$  in world frame  $\mathcal{F}_w$ ; a *Local Mapping* thread which builds on the map and the keyframes and frequently performs local bundle adjustments to update the positions of map points and camera poses at the keyframes; a *Loop Closing* thread which detects loops in the trajectory and propagates a correction through the trajectory poses and the map. In addition, we implemented the following improvements:

<sup>7</sup> [https://github.com/ristofer/visual\\_localization/](https://github.com/ristofer/visual_localization/)

<sup>8</sup> [http://wiki.ros.org/move\\_base](http://wiki.ros.org/move_base)

**Tracking modifications** We changed the Tracking thread to process not only images but also odometry measurements, obtained directly from ROS. Odometry measurements are computed within the Pepper’s internal software and published in ROS through NAOqi wrappers with respect to the *odom* frame, shown as  $\mathcal{F}_o$  in Figure 2. Our ROS-compatible wrapper for ORB-SLAM subscribes the *tf* topic and images, and requests an odometry measurement every time a new image is received, obtaining a synchronized pair image-odometry. Later, every time a new keyframe is created after a successful camera tracking, the odometry information is also included in the keyframe.

In addition, since the original behavior of ORB-SLAM is to stop providing camera poses when camera tracking fails, and wait until a relocalization, which is not a desirable strategy while navigating<sup>9</sup>, we set the camera estimation equal to the odometry prediction. This ensures a continuous camera pose hypothesis for planning tasks but requires that the metric scale is initialized.

**Metric scale initialization** We did not utilize any general system initialization solution as in [9] but preferred a multi-step approach as in [12]. We first wait until the pure visual SLAM system is initialized and the unscaled map built, to then compute the scale from the odometry information between keyframes. By comparing the relative translations between keyframes as predicted by ORB-SLAM  $\Delta\mathbf{p}_o(i-1, i)$  and the odometry  $\Delta\mathbf{p}_w(i-1, i)$ , the scale can be retrieved and the map and keyframe poses can be updated by the method of Horn [6] (Eq. (2)). However, the initial map is subject to major change in the early stages of the mapping. Therefor the scale correction is done after a fixed number of  $N$  keyframes have been created, ensuring a satisfactory converged map and thus a reasonably reliable scale correction. The success of this strategy only depends on the environment’s size and the motion performed by the robot; an additional discussion is given in Section 5.

$$s = \frac{\sqrt{\sum_i^N \|\Delta\mathbf{p}_o(i-1, i)\|^2}}{\sqrt{\sum_i^N \|\Delta\mathbf{p}_w(i-1, i)\|^2}} \quad (2)$$

After the scale update, a Global Bundle Adjustment (Global BA) is performed to guarantee an optimal map reconstruction.

**Local Mapping** Every time a new keyframe is created, *Local Mapping* performs an optimization in a subset of the complete trajectory updating both the poses and the map -the so-called *local window*. The parts of the trajectory to be optimized are keyframes in the neighborhood of the last added keyframe, and also map points being observed by those; the neighbors are selected by the so-called *covisibility graph* [10]. This operation on the local window ensures an efficient optimization process even in large maps.

<sup>9</sup> Unless high-level behaviors to detect failures are considered

Since the initialization procedure makes the current trajectory and map *metrically consistent*, it is possible to fuse the visual information with wheel odometry information to avoid drift. This is done by adding odometry factors or *constraints* between keyframes. In order to do so, the odometry measurement is mapped from the odometry frame  $\vec{\mathcal{F}}_o$  to the camera frame  $\vec{\mathcal{F}}_c$  by using Pepper's forward kinematics. Hence, we compute the difference between the odometry measurements  $i-1$  and  $i$ ,  ${}_c\mathbf{T}_{C_{i-1},C_i}$ , between all the keyframes in the local window, which hopefully match the difference between the keyframes' pose,  $(\mathbf{T}_{C_{i-1}W}\mathbf{T}_{C_iW})$ . The error between the odometry's and ORB-SLAM's differences are defined in the minimal representation of the pose, i.e. 6-dimensional, which is achieved by using the *logarithm map* of SE(3):

$$\varepsilon_{odo} = \text{Log}_{\text{SE}(3)} \left( {}_c\mathbf{T}_{C_{i-1},C_i}^{-1} \mathbf{T}_{C_{i-1}W} \mathbf{T}_{C_iW}^{-1} \right). \quad (3)$$

This residual is defined for every pair of keyframes within the local window; additionally, keyframes with neighbors which are not in the local window, are also added as fixed nodes in the optimization. The corresponding optimization problem that minimizes visual error terms  $\varepsilon_{vis}$  (as defined in [10]) and odometry terms  $\varepsilon_{odo}$  (Eq. (3)), is:

$$\mathcal{X}^* = \arg \min_{\mathcal{X}} \sum_{(i,k)} \|\varepsilon_{vis}\|_{\mathbf{\Omega}_{vis}}^2 + \sum_{(i-1,i)} \|\varepsilon_{odo}\|_{\mathbf{\Omega}_{odo}}^2 \quad (4)$$

The optimization problem in equation (4) is solved with the graph optimization framework *g2o* [8] using fixed information matrices  $\mathbf{\Omega}_{odo}, \mathbf{\Omega}_{vis}$ . The resulting keyframe poses and map points are then updated, and the *Local Mapping* thread awaits until a new keyframe is added from *Tracking*.

**Localization mode and map reuse** ORB-SLAM provides the option to localize in a previously built map, disabling the SLAM capabilities. This localization can run in a single thread, hence requiring a fraction of the computational requirements compared to the full ORB-SLAM system. Nevertheless, in order to perform localization-only, it is required a map that was built in the same session.

Since this is not generally the case, we use map saving capabilities (taken from a fork of ORB-SLAM<sup>10</sup>) and implemented a different behavior for the system when it is launch with a pre-built map, that first tries to relocalize and then continue mapping incrementally. These minor changes allowed us to build maps even during different sessions once the relocalization is successful.

## 4.2 Navigation

For the navigation part, we assume the ORB-SLAM's map was already built, so we can rely on the localization mode. The pose estimation performed by ORB-SLAM is 6-dimensional since it considers a camera freely moving in the space,

<sup>10</sup> [https://github.com/Alkaid-Benetnash/ORB\\_SLAM2](https://github.com/Alkaid-Benetnash/ORB_SLAM2)

which would be an overkill to perform planning with Pepper. In order to use ORB-SLAM's estimates within a planar navigation framework, we developed the `visual_localization` node, which takes the estimated position of the camera with respect to the ORB-SLAM world frame and computes a transformation between the ROS standard *map* and *odom* frames. This transformation represents the Pepper position in the ORB-SLAM map based on the estimated pose of the Pepper's camera and its kinematic information.

The `move_base` package is used to navigate. Our localization system basically replaces the `amcl`<sup>11</sup> package in the ROS Navigation Stack. The `move_base` package uses the pose estimate provided by the localization system and Pepper's laser readings to compute the cost map necessary for planning. Thus, lasers are not used for localization, but for obstacle detection and path planning.

## 5 Experiments and results

### 5.1 Experimental setup

We considered two real environments of the Faculty of Physical and Mathematical Sciences of Universidad de Chile to test our system: *Mechatronics Laboratory* and *School Building South Hall*. The chosen places were different in size, furniture, and visual features complexity, being the latter of paramount importance for the visual SLAM system.

- The Mechatronics Laboratory (Figure 4a) is a  $10 \times 9m^2$  space. The main furniture are rolling chairs and work tables. It is a feature-rich space comparable to the RoboCup Arena; however, it has various windows that enable the pass of natural light.
- The School Building South Hall (Figure 4b) has an area of  $16 \times 27.5m^2$ . It is an open space with pillars and doors, but generally feature-less, making it the most challenging environment for our system.

To have a ground truth reference, a Google Tango Tablet is used (Figure 4c).



**Fig. 4.** *Left:* Mechatronics Lab. *Center:* South Hall. *Right:* Pepper with Google Tango Tablet attached for ground truth measurements.

<sup>11</sup> <http://wiki.ros.org/amcl>



## 5.2 Experiments

**Mapping** The first experiment considered a localization and mapping task; this was performed in both the Mechatronics Lab (Fig. 4a) and South Hall (Fig. 4b). We remote controlled the robot to build a three dimensional map to be later used for localization. Table 1 compare different mapping results through the *Absolute Trajectory Error* [14], a metric that calculates the root mean square error *RMSE* defined as  $\left(\frac{1}{N} \sum_i^N \|p_{e_i} - p_{gt_i}\|^2\right)^{1/2}$  between the localization estimate  $p_{e_i}$  and the ground truth  $p_{gt_i}$  through all the time indices.

**Table 1.** Absolute Trajectory Error (ATE) in meters, for each place and axis. A mapping experiment was performed in the Mechatronics Laboratory and in the South Hall. The estimated trajectory and the ground truth was used to calculate the ATE.

Place	ATE X [m]	ATE Y [m]	ATE Z [m]
Mechatronics Laboratory	0.270	0.249	0.080
South Hall	0.619	0.849	0.390

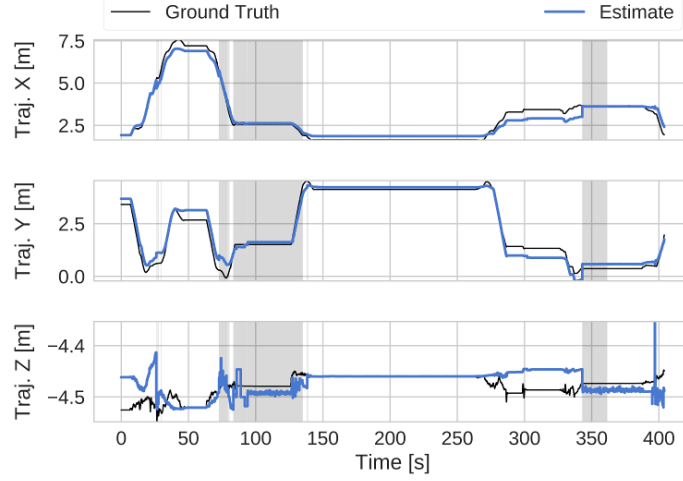
During all the experiments we noticed that the robot must move smoothly and preferably sideways in order to triangulate the initial map; pure rotational factors must be avoided despite the offset between the head camera and the base’s axis of rotation. The initial displacement is primordial to recover a reliable scale factor as well. However, this also depends on a parameter that sets the number of keyframes to wait until the scale is recovered with Eq. 2, which is set empirically.

Regarding mapping, as is expected from a feature-based visual SLAM system, the number of points and quality of the map increases in feature-rich environments. In addition, compared to LIDAR mapping, visual mapping requires significantly more time. This because map creation depends on the field-of-view (FOV) of the camera, which is very narrow in Pepper, requiring to map the same place from multiple views in order make it useful for robust localization. LIDAR does not suffer from this issue since localization is performed by point cloud alignment rather than feature matching. However, feature matching has the advantage of providing instantaneous relocalization when the robot is lost since places are uniquely defined by a bag-of-words representation [10].

**Localization and Navigation** We performed a second experiment to test the localization and navigation in a known place, i.e., with a pre-built map. This was also executed in the Mechatronics Laboratory and South Hall.

We commanded the robot to navigate without operator help to a relative point with respect to its initial pose, which exploited the localization capabilities of our system in a known environment. Localization results are shown in Figures 5 and 6.

Our experiments show the performance of the system, which uses both visual localization and odometry fusion (highlighted in gray) and odometry-only localization when the robot is lost (in white). In the navigation experiment in the



**Fig. 5.** Navigation test on the Mechatronics Laboratory. The estimate of the localization system is compared to ground truth. Grey areas in the graph indicated when the robot is correctly localized with ORB-SLAM. When the robot is not localized, an odometry estimate is used

Mechatronics Lab, showed in Figure 5, Pepper correctly navigated through the test. However, between seconds 275 and 350 there exists a considerable drift between the ground truth and the localization estimate. These problems can result in reaching an erroneous navigation goal or even collide if no safety procedures are considered. We believe that a cause of this issue was the lack of viewpoints during the mapping step, as mentioned in the previous experiment.

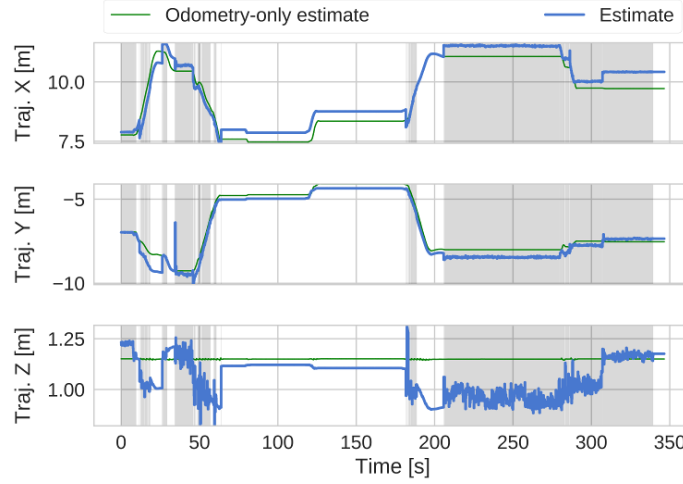
Regarding South Hall experiments, the multiple discontinuities in the localization estimate (Figure 6, Z axis) made navigation unfeasible. This was caused by the large distance between the robot and the landmarks in this environment, which was not the case of the Mechatronics Lab. Since visual SLAM systems are based on optimization and Pepper’s FOV is narrow, it is more difficult to correctly estimate the pose because the triangulation uncertainty is higher; this is a known problem in visual systems [4].

However, in both experiments we noticed that localization is robust to changes in the environment, like a change in furniture position and, if the map is correctly built, there is minimal (if not zero) accumulation of drift error.

### 5.3 Discussion

The previous experiments evidence several advantages but also challenges of our proposed system. We summarize them as follows:

*Advantages* Our localization system is not affected by Peppers’ LIDARs short range, which is one of the main limitations of it in RoboCup environments. Since we used the map created visually, the robot is able to localize with a single look



**Fig. 6.** Navigation test on the South Hall. The robot tries to navigate but the localization system does not work correctly.

by exploiting the relocalization capabilities of ORB-SLAM. Our scale recovery solution also allowed us to perform metric mapping despite using a single camera. In addition, since the motion estimation is based on features, it is robust to partial occlusions, and odometry is used when no visual features are tracked. All these advantages demonstrate that it is possible for our system to localize the robot in RoboCup@Home arenas successfully.

*Challenges* Despite the previous advantages, we cannot avoid to mention some challenges and difficulties we noticed during our experiments. The first one relates to illumination changes, which can deteriorate hugely visual tracking. If we set the camera exposure to automatic we can deal with dynamic lighting, but the system is more susceptible to motion blur, which is still an issue despite the robot performs planar motion; the main cause of this is joint backlash. If the environment has non-variable illumination, we recommend to fix the camera exposure to diminish those problems. The second challenge regards glossy surfaces, which produce fake landmarks because of reflections. Despite ORB-SLAM is able to deal with outliers that do not match the predicted motion, it is still an open challenge in our opinion. Finally, localization turns noisier when the landmarks are far away, which is caused by the optimization procedures and the point’s triangulation uncertainty.

## 6 Conclusions

In this work, we presented a localization system for a Pepper robot based on a visual SLAM system. Our solution, built upon ORB-SLAM, was focused on developing a self-localization system able to deal with large environments despite the LIDARs’ short range. In order to do so, we presented an approach

that fused visual and wheel odometry information. We tested the system in two real environments, showing the feasibility performing SLAM and navigation with our system with the current Pepper sensors, despite displaying some issues such as weakness to illumination changes, ambiguities to glassy surfaces and far landmarks.

Nowadays we are working towards an on-board implementation of the self-localization system on Pepper, which will allow us to perform a more exhaustive evaluation and comparison with other sensors such as lasers. In the future, we would like to improve robustness to illumination changes and reducing the noisy behavior in large environments.

## Acknowledgments

This work was partially funded by the FONDECYT 1161500 project.

## References

1. Cadena, C., Carlone, L., Carrillo, H., Latif, Y., Scaramuzza, D., Neira, J., Reid, I., Leonard, J.J.: Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age. *IEEE Trans. Robot.* 32(6), 1309–1332 (2016)
2. Cheng, M., Chen, X., Tang, K., Wu, F., Kupcsik, A., Iocchi, L., Chen, Y., Hsu, D.: Synthetical benchmarking of service robots: A first effort on domestic mobile platforms. In: *Lecture Notes in Computer Science*. vol. 9513, pp. 377–388 (2015)
3. Engel, J., Schöps, T., Cremers, D.: LSD-SLAM: Large-Scale Direct Monocular SLAM. In: *Eccv*, pp. 834–849 (2014)
4. Hartley, R., Zisserman, A.: *Multiple View Geometry in Computer Vision*. Cambridge University Press (2004)
5. Hertzberg, C., Wagner, R., Frese, U., Schröder, L.: Integrating generic sensor fusion algorithms with sound state representations through encapsulation of manifolds. *Information Fusion* 14(1), 57–77 (2013)
6. Horn, B.K.P.: Closed-form solution of absolute orientation using unit quaternions. *Journal of the Optical Society of America A* 4(4), 629 (1987)
7. Iocchi, L., Holz, D., Ruiz-Del-Solar, J., Sugiura, K., Van Der Zant, T.: RoboCup@Home: Analysis and results of evolving competitions for domestic and service robots. *Artificial Intelligence* 229, 258–281 (2015)
8. Kümmerle, R., Grisetti, G., Strasdat, H., Konolige, K., Burgard, W.: G2o: A general framework for graph optimization. *IEEE ICRA* pp. 3607–3613 (2011)
9. Martinelli, A.: Closed-form solution of visual-inertial structure from motion. *International Journal of Computer Vision* 106(2), 138–152 (2014)
10. Mur-Artal, R., Montiel, J.M., Tardos, J.D.: ORB-SLAM: A Versatile and Accurate Monocular SLAM System. *IEEE Transactions on Robotics* 31(5), 1147–1163 (2015)
11. Mur-Artal, R., Tardos, J.D.: ORB-SLAM2: An Open-Source SLAM System for Monocular, Stereo, and RGB-D Cameras. *IEEE Trans. Robot.* 33(5), 1255–1262 (2017)
12. Mur-Artal, R., Tardos, J.D.: Visual-Inertial Monocular SLAM With Map Reuse. *IEEE Robotics and Automation Letters* 2(2), 796–803 (apr 2017)
13. Strasdat, H., Montiel, J.M.M., Davison, A.J.: *Visual SLAM: Why Filter?* (2012)
14. Sturm, J., Engelhard, N., Endres, F., Burgard, W., Cremers, D.: A benchmark for the evaluation of RGB-D SLAM systems. *IEEE ICRA* pp. 573–580 (2012)